# Load The Dataset (Week 2)

In [1]:
```python
import pandas as pd

#ingest data
df=pd.read_csv('C:/Users/zhumh/Downloads/hotel_booking.csv.zip')
df.head()
```

Out[1]:

| | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_week_number | arrival |
|---|---|---|---|---|---|---|---|
| 0 | Resort Hotel | 0 | 342 | 2015 | July | 27 | |
| 1 | Resort Hotel | 0 | 737 | 2015 | July | 27 | |
| 2 | Resort Hotel | 0 | 7 | 2015 | July | 27 | |
| 3 | Resort Hotel | 0 | 13 | 2015 | July | 27 | |
| 4 | Resort Hotel | 0 | 14 | 2015 | July | 27 | |

5 rows × 36 columns

◀ ▬▬▬▬▬▬▬▬ ▶

In [2]:
```python
#basic information of dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119390 entries, 0 to 119389
Data columns (total 36 columns):
 #   Column                       Non-Null Count   Dtype
---  ------                       --------------   -----
 0   hotel                        119390 non-null  object
 1   is_canceled                  119390 non-null  int64
 2   lead_time                    119390 non-null  int64
 3   arrival_date_year            119390 non-null  int64
 4   arrival_date_month           119390 non-null  object
 5   arrival_date_week_number     119390 non-null  int64
 6   arrival_date_day_of_month    119390 non-null  int64
 7   stays_in_weekend_nights      119390 non-null  int64
 8   stays_in_week_nights         119390 non-null  int64
 9   adults                       119390 non-null  int64
 10  children                     119386 non-null  float64
 11  babies                       119390 non-null  int64
 12  meal                         119390 non-null  object
```

```
 13   country                      118902 non-null   object
 14   market_segment               119390 non-null   object
 15   distribution_channel         119390 non-null   object
 16   is_repeated_guest            119390 non-null   int64
 17   previous_cancellations       119390 non-null   int64
 18   previous_bookings_not_canceled  119390 non-null   int64
 19   reserved_room_type           119390 non-null   object
 20   assigned_room_type           119390 non-null   object
 21   booking_changes              119390 non-null   int64
 22   deposit_type                 119390 non-null   object
 23   agent                        103050 non-null   float64
 24   company                      6797 non-null     float64
 25   days_in_waiting_list         119390 non-null   int64
 26   customer_type                119390 non-null   object
 27   adr                          119390 non-null   float64
 28   required_car_parking_spaces  119390 non-null   int64
 29   total_of_special_requests    119390 non-null   int64
 30   reservation_status           119390 non-null   object
 31   reservation_status_date      119390 non-null   object
 32   name                         119390 non-null   object
 33   email                        119390 non-null   object
 34   phone-number                 119390 non-null   object
 35   credit_card                  119390 non-null   object
dtypes: float64(4), int64(16), object(16)
memory usage: 32.8+ MB
```

In [3]:
```python
df.isnull().mean()
```

Out[3]:
```
hotel                            0.000000
is_canceled                      0.000000
lead_time                        0.000000
arrival_date_year                0.000000
arrival_date_month               0.000000
arrival_date_week_number         0.000000
arrival_date_day_of_month        0.000000
stays_in_weekend_nights          0.000000
stays_in_week_nights             0.000000
adults                           0.000000
children                         0.000034
babies                           0.000000
meal                             0.000000
country                          0.004087
market_segment                   0.000000
distribution_channel             0.000000
is_repeated_guest                0.000000
previous_cancellations           0.000000
previous_bookings_not_canceled   0.000000
reserved_room_type               0.000000
assigned_room_type               0.000000
booking_changes                  0.000000
deposit_type                     0.000000
agent                            0.136862
company                          0.943069
days_in_waiting_list             0.000000
customer_type                    0.000000
adr                              0.000000
required_car_parking_spaces      0.000000
total_of_special_requests        0.000000
reservation_status               0.000000
```

```
reservation_status_date         0.000000
name                            0.000000
email                           0.000000
phone-number                    0.000000
credit_card                     0.000000
dtype: float64
```

In [4]:
```python
# transpose the resulting DataFrame
df.describe([0.01,0.05,0.1,0.25,0.5,0.75,0.99]).T
```

Out[4]:

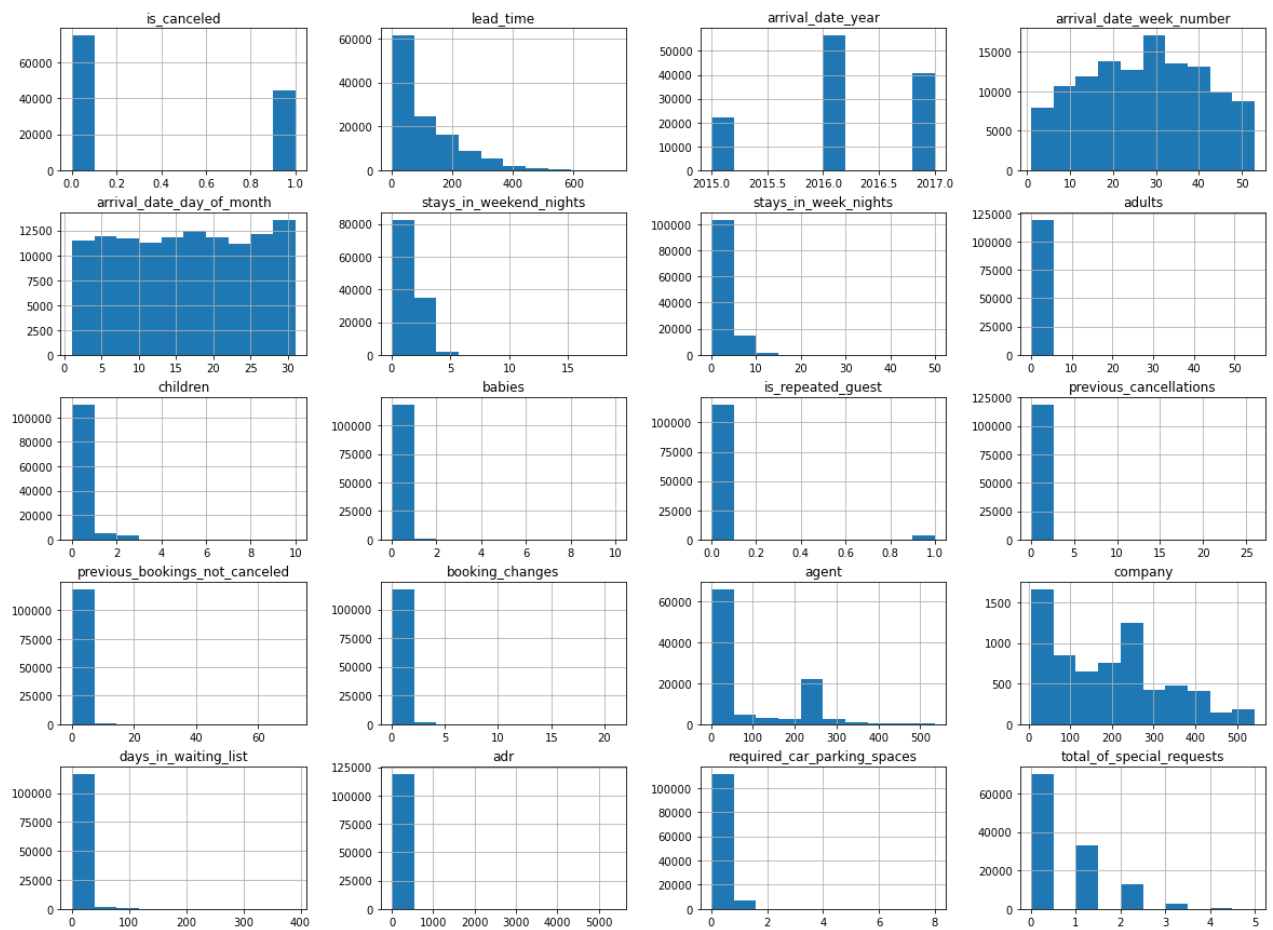| | count | mean | std | min | 1% | 5% | 10% | 25 |
|---|---|---|---|---|---|---|---|---|
| is_canceled | 119390.0 | 0.370416 | 0.482918 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| lead_time | 119390.0 | 104.011416 | 106.863097 | 0.00 | 0.0 | 0.0 | 3.0 | 18 |
| arrival_date_year | 119390.0 | 2016.156554 | 0.707476 | 2015.00 | 2015.0 | 2015.0 | 2015.0 | 2016 |
| arrival_date_week_number | 119390.0 | 27.165173 | 13.605138 | 1.00 | 2.0 | 5.0 | 8.0 | 16 |
| arrival_date_day_of_month | 119390.0 | 15.798241 | 8.780829 | 1.00 | 1.0 | 2.0 | 4.0 | 8 |
| stays_in_weekend_nights | 119390.0 | 0.927599 | 0.998613 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| stays_in_week_nights | 119390.0 | 2.500302 | 1.908286 | 0.00 | 0.0 | 0.0 | 1.0 | 1 |
| adults | 119390.0 | 1.856403 | 0.579261 | 0.00 | 1.0 | 1.0 | 1.0 | 2 |
| children | 119386.0 | 0.103890 | 0.398561 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| babies | 119390.0 | 0.007949 | 0.097436 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| is_repeated_guest | 119390.0 | 0.031912 | 0.175767 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| previous_cancellations | 119390.0 | 0.087118 | 0.844336 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| previous_bookings_not_canceled | 119390.0 | 0.137097 | 1.497437 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| booking_changes | 119390.0 | 0.221124 | 0.652306 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| agent | 103050.0 | 86.693382 | 110.774548 | 1.00 | 1.0 | 1.0 | 6.0 | 9 |
| company | 6797.0 | 189.266735 | 131.655015 | 6.00 | 16.0 | 40.0 | 40.0 | 62 |
| days_in_waiting_list | 119390.0 | 2.321149 | 17.594721 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| adr | 119390.0 | 101.831122 | 50.535790 | -6.38 | 0.0 | 38.4 | 50.0 | 69 |
| required_car_parking_spaces | 119390.0 | 0.062518 | 0.245291 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |
| total_of_special_requests | 119390.0 | 0.571363 | 0.792798 | 0.00 | 0.0 | 0.0 | 0.0 | 0 |

In [5]:
```python
import matplotlib.pyplot as plt

# generate histograms for all the columns
df.hist(figsize=(20,15))
plt.show()
```
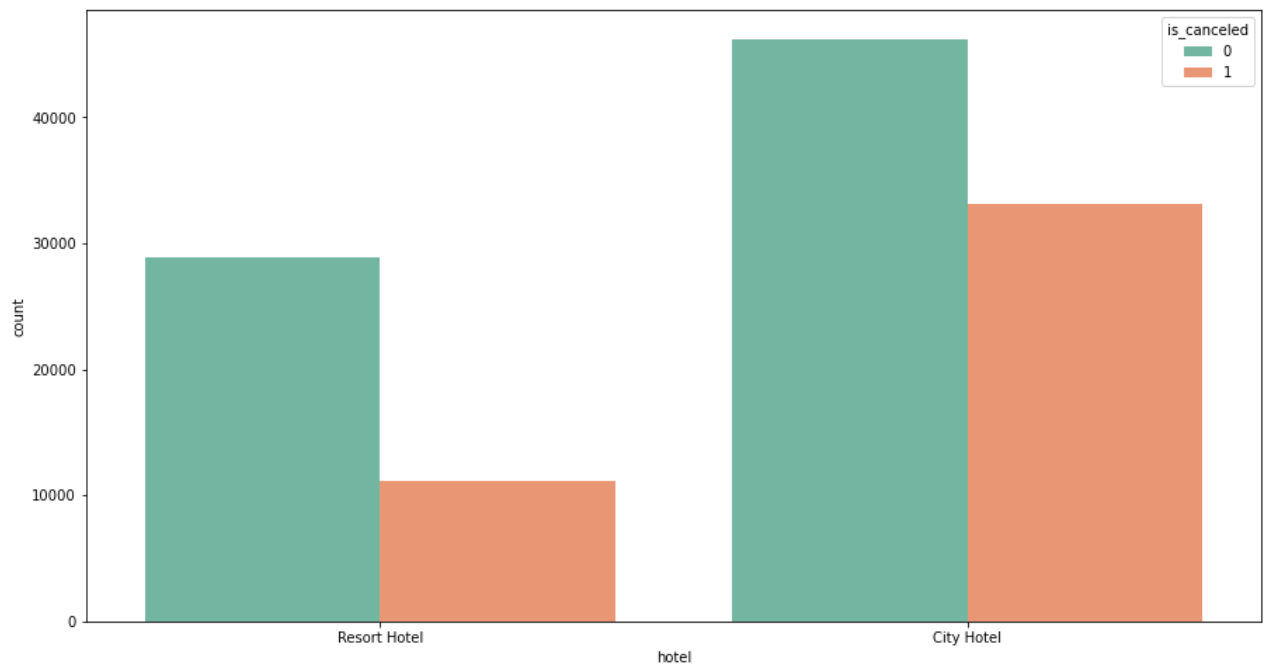
# Explore The Data

1. Hotel bookings and cancellations

In [6]:
```python
import seaborn as sns
plt.figure(figsize=(15,8))
sns.countplot(x='hotel'
              ,data=df
              ,hue='is_canceled'
              ,palette=sns.color_palette('Set2',2)
              )
```

Out[6]:  `<AxesSubplot:xlabel='hotel', ylabel='count'>`

In [7]:
```python
hotel_cancel=(df.loc[df['is_canceled']==1]['hotel'].value_counts()/df['hotel'].value_co
print('Hotel cancellations'.center(20),hotel_cancel,sep='\n')
```

```
Hotel cancellations
City Hotel      0.417270
Resort Hotel    0.277634
Name: hotel, dtype: float64
```

City Hotel's booking volume and cancellation volume are both higher than Resort Hotel's, but Resort Hotel's cancellation rate is 27.8%, while City Hotel's cancellation rate reaches 41.7%.

1. Hotel bookings by month

In [8]:
```python
city_hotel=df[(df['hotel']=='City Hotel') & (df['is_canceled']==0)]
resort_hotel=df[(df['hotel']=='Resort Hotel') & (df['is_canceled']==0)]
for i in [city_hotel,resort_hotel]:
    i.index=range(i.shape[0])

city_month=city_hotel['arrival_date_month'].value_counts()
resort_month=resort_hotel['arrival_date_month'].value_counts()
name=resort_month.index
x=list(range(len(city_month.index)))
y=city_month.values
x1=[i+0.3 for i in x]
y1=resort_month.values
width=0.3
plt.figure(figsize=(15,8),dpi=80)
plt.plot(x,y,label='City Hotel',color='lightsalmon')
plt.plot(x1,y1,label='Resort Hotel',color='lightseagreen')
plt.xticks(x,name)
plt.legend()
plt.xlabel('Month')
plt.ylabel('Count')
plt.title('Month Book')
for x,y in zip(x,y):
    plt.text(x,y+0.1,'%d' % y,ha = 'center',va = 'bottom')
```

```
for x,y in zip(x1,y1):
    plt.text(x,y+0.1,'%d' % y,ha = 'center',va = 'bottom')
```
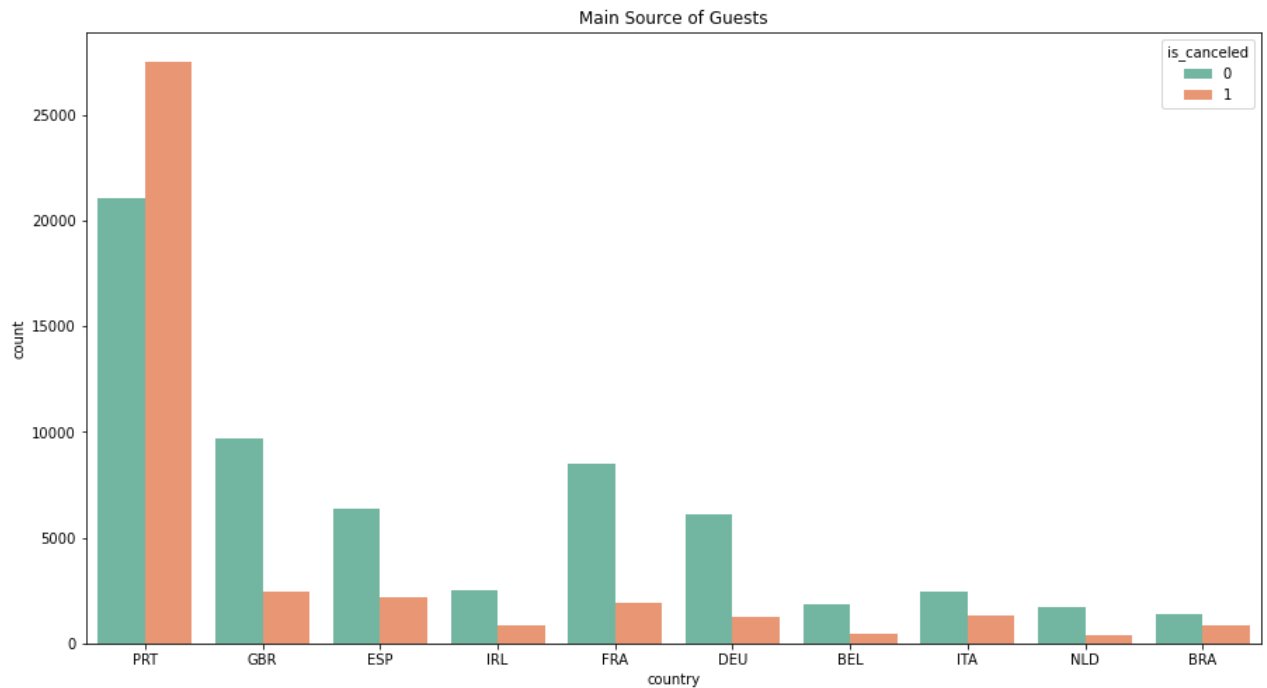


Peak booking months are August and July.

1. Customer origin and booking cancellation rate

In [9]:
```
country_book=df['country'].value_counts()[:10]
country_cancel=df[(df.country.isin (country_book.index)) & (df.is_canceled==1)]['countr
plt.figure(figsize=(15,8))
sns.countplot(x='country'
             ,data=df[df.country.isin (country_book.index)]
             ,hue='is_canceled'
             ,palette=sns.color_palette('Set2',2)
             )
plt.title('Main Source of Guests')
```

Out[9]:   Text(0.5, 1.0, 'Main Source of Guests')

```
In [10]:   country_cancel_rate=(country_cancel/country_book).sort_values(ascending=False)
           print('Customer cancellation rates by country'.center(10),country_cancel_rate,sep='\n')
```

```
Customer cancellation rates by country
PRT     0.566351
BRA     0.373201
ITA     0.353956
ESP     0.254085
IRL     0.246519
BEL     0.202391
GBR     0.202243
FRA     0.185694
NLD     0.183935
DEU     0.167147
Name: country, dtype: float64
```
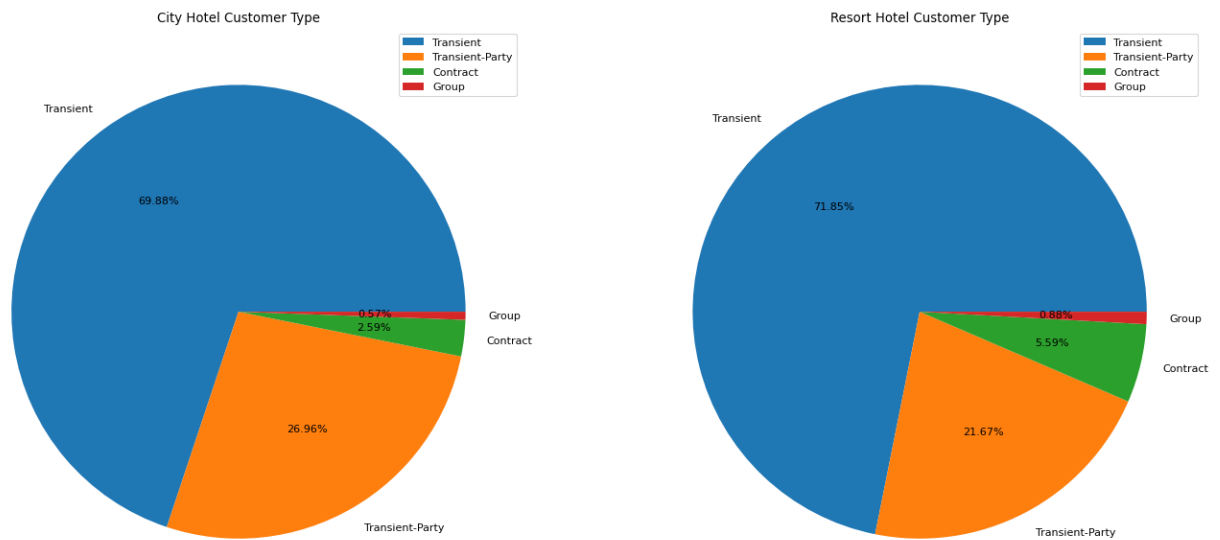
The peak season for both Resort hotel and City hotel is July and August in summer, and the main sources of tourists are European countries. This is in line with the characteristics of European tourists who prefer summer travel. It is necessary to focus on countries with high cancellation rates such as Portugal (PRT) and the United Kingdom (BRT). Main source of customers.

1. Customer type

```
In [11]:   city_customer=city_hotel.customer_type.value_counts()
           resort_customer=resort_hotel.customer_type.value_counts()
           plt.figure(figsize=(21,12),dpi=80)
           plt.subplot(1,2,1)
           plt.pie(city_customer,labels=city_customer.index,autopct='%.2f%%')
           plt.legend(loc=1)
           plt.title('City Hotel Customer Type')
           plt.subplot(1,2,2)
           plt.pie(resort_customer,labels=resort_customer.index,autopct='%.2f%%')
           plt.title('Resort Hotel Customer Type')
```
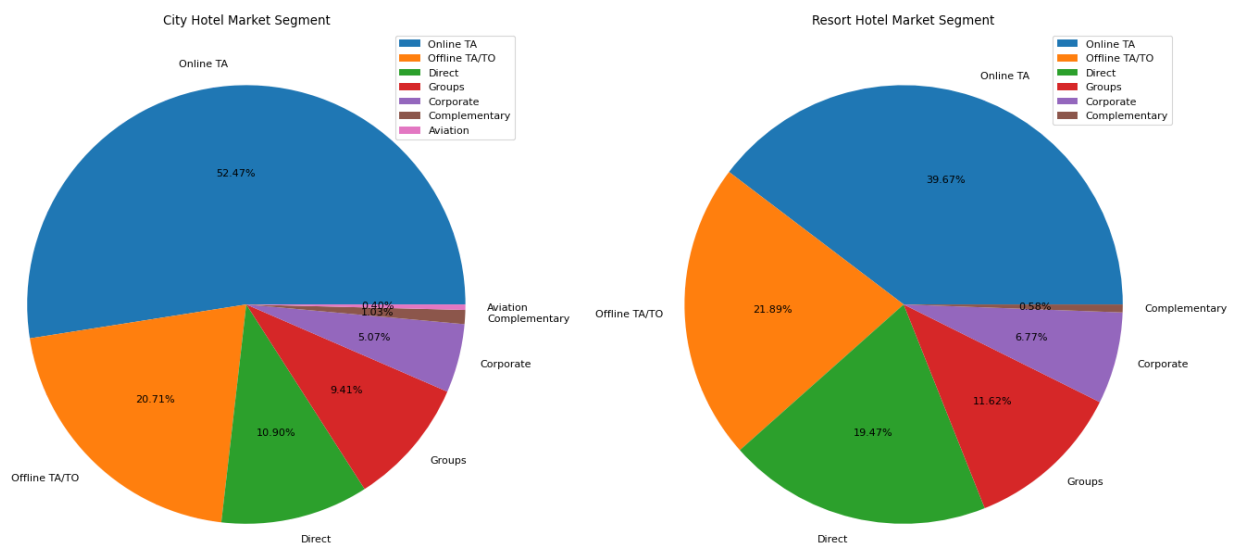
```
plt.legend()
plt.show()
```



The main customer type of the hotel is transient travelers, accounting for about 70%.

1. Hotel booking method

In [12]:
```
city_segment=city_hotel.market_segment.value_counts()
resort_segment=resort_hotel.market_segment.value_counts()
plt.figure(figsize=(21,12),dpi=80)
plt.subplot(1,2,1)
plt.pie(city_segment,labels=city_segment.index,autopct='%.2f%%')
plt.legend()
plt.title('City Hotel Market Segment')
plt.subplot(1,2,2)
plt.pie(resort_segment,labels=resort_segment.index,autopct='%.2f%%')
plt.title('Resort Hotel Market Segment')
plt.legend()
plt.show()
```

The customers of the two hotels mainly come from online travel agencies, which account for even more than 50% of the City Hotel; offline travel agencies come next, accounting for about 20%.

1. Average daily expenses of various types of passengers

In [13]:
```python
plt.figure(figsize=(15,8))
sns.boxplot(x='customer_type'
            ,y='adr'
            ,hue='hotel'
            ,data=df[df.is_canceled==0]
            ,palette=sns.color_palette('Set2',2)
            )
plt.title('Average Daily Rate of Different Customer Type')
```

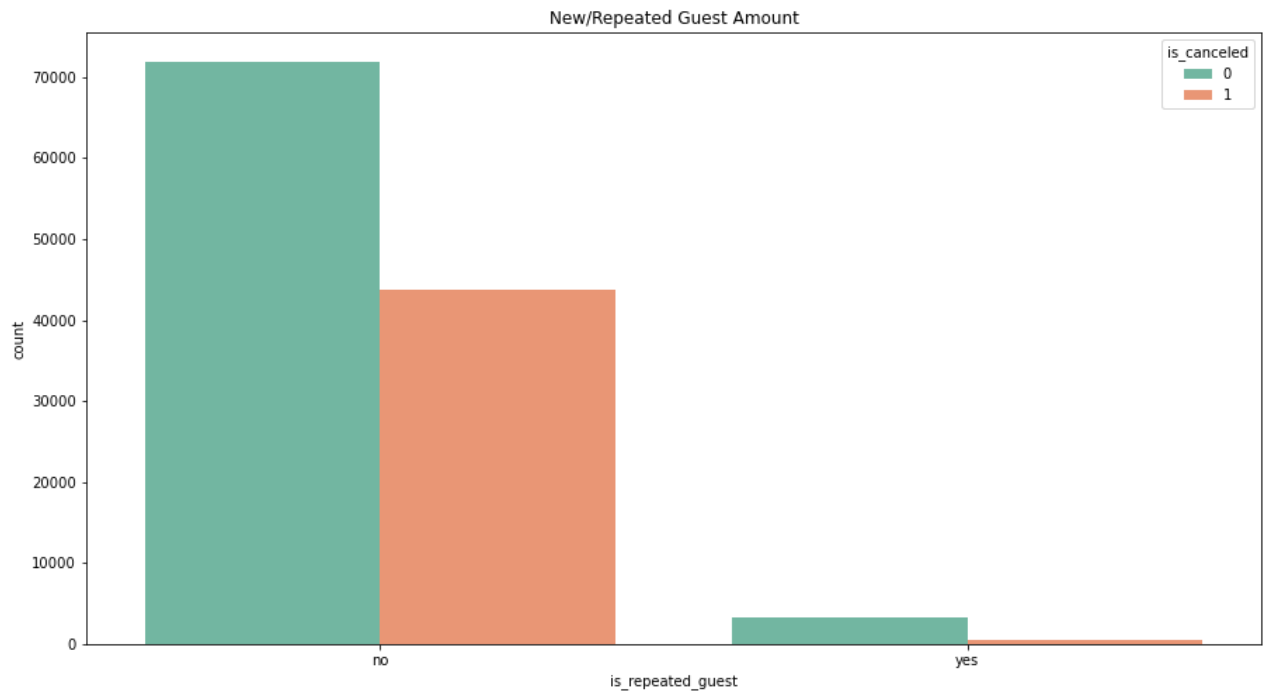Out[13]:   Text(0.5, 1.0, 'Average Daily Rate of Different Customer Type')



The average daily expenditure of all types of customers of City Hotel is higher than that of Resort Hotel; among the four types of customers, the consumption of individual travelers (Transient) is the highest and that of group travelers (Group) is the lowest.

7.Number of new and old customers and cancellation rate

In [14]:
```python
plt.figure(figsize=(15,8))
sns.countplot(x='is_repeated_guest'
              ,data=df
              ,hue='is_canceled'
              ,palette=sns.color_palette('Set2',2)
              )
plt.title('New/Repeated Guest Amount')
plt.xticks(range(2),['no','yes'])
```

Out[14]:
```
([<matplotlib.axis.XTick at 0x235beb346d0>,
  <matplotlib.axis.XTick at 0x235beb346a0>],
 [Text(0, 0, 'no'), Text(1, 0, 'yes')])
```

New/Repeated Guest Amount



In [15]:
```
guest_cancel=(df.loc[df['is_canceled']==1]['is_repeated_guest'].value_counts()/df['is_r
guest_cancel.index=['New Guest', 'Repeated Guest']
print('Cancellation rate for new and old customers'.center(15),guest_cancel,sep='\n')
```

```
Cancellation rate for new and old customers
New Guest          0.377851
Repeated Guest     0.144882
Name: is_repeated_guest, dtype: float64
```

The cancellation rate for regular customers was 14.4%, while the cancellation rate for new customers reached 37.8%, which was 24 percentage points higher than that for regular customers.

1. Deposit method and reservation cancellation rate

In [16]:
```
print('Three deposit methods for booking quantity'.center(15),df['deposit_type'].value_
```

```
Three deposit methods for booking quantity
No Deposit      104641
Non Refund       14587
Refundable         162
Name: deposit_type, dtype: int64
```
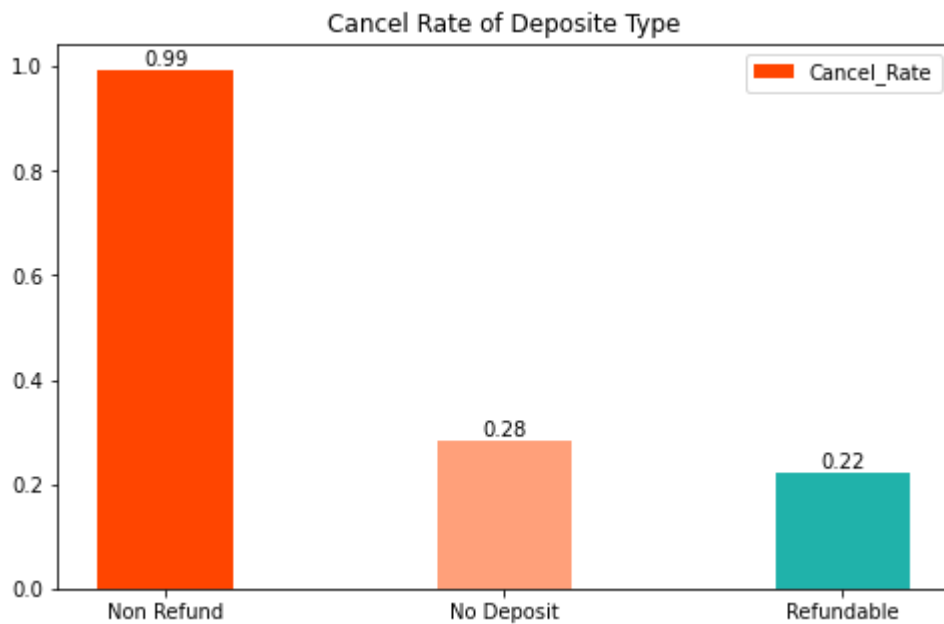
In [17]:
```
deposit_cancel=(df.loc[df['is_canceled']==1]['deposit_type'].value_counts()/df['deposit
plt.figure(figsize=(8,5))
x=range(len(deposit_cancel.index))
y=deposit_cancel.values
plt.bar(x,y,label='Cancel_Rate',color=['orangered','lightsalmon','lightseagreen'],width
plt.xticks(x,deposit_cancel.index)
plt.legend()
plt.title('Cancel Rate of Deposite Type')
for x,y in zip(x,y):
    plt.text(x,y,'%.2f' % y,ha = 'center',va = 'bottom')
```
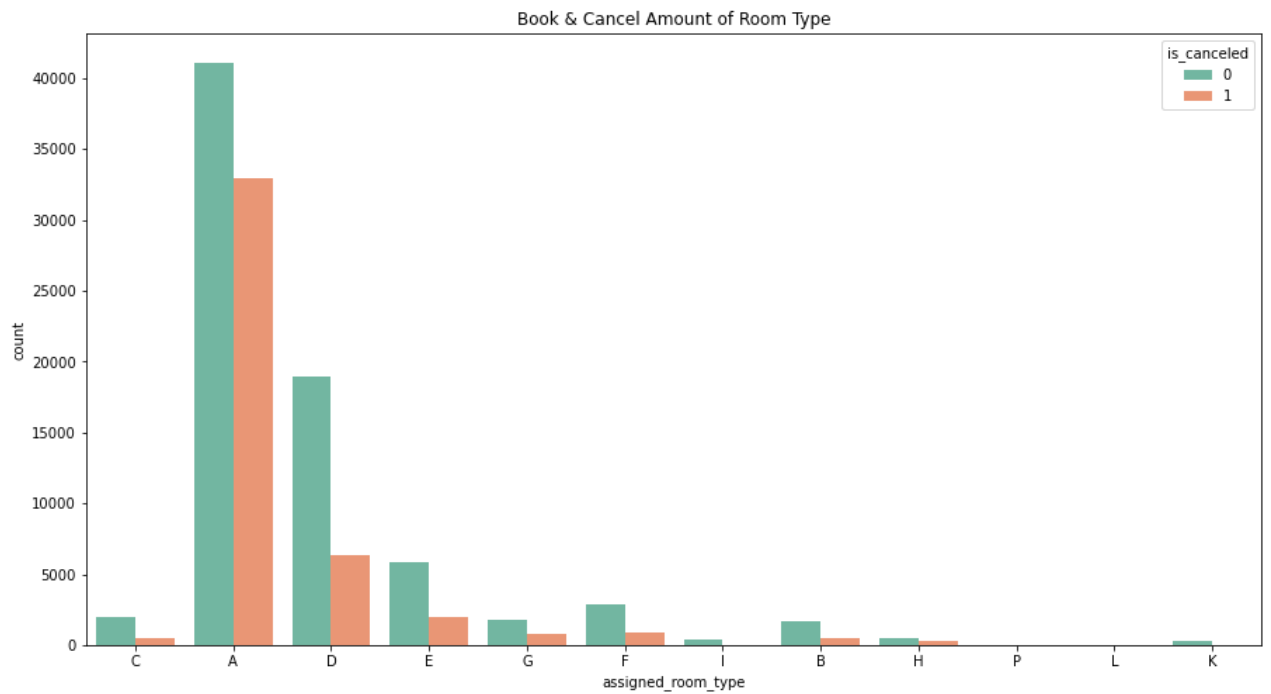
Cancel Rate of Deposite Type



'No Deposit' is the method with the highest number of bookings and has a low cancellation rate, while the cancellation rate of non-refundable type is as high as 99%. This type of deposit method can be reduced to reduce Customer cancellation rate.

1. Room type and cancellation volume

In [18]:
```python
plt.figure(figsize=(15,8))
sns.countplot(x='assigned_room_type'
             ,data=df
             ,hue='is_canceled'
             ,palette=sns.color_palette('Set2',2)
             )
plt.title('Book & Cancel Amount of Room Type')
```

Out[18]:    Text(0.5, 1.0, 'Book & Cancel Amount of Room Type')



In [19]:
```
room_cancel=df.loc[df['is_canceled']==1]['assigned_room_type'].value_counts()[:7]/df['a
print('Cancellation rates for different room types'.center(5),room_cancel.sort_values(a
```

```
Cancellation rates for different room types
A    0.444925
G    0.305523
E    0.252114
D    0.251244
F    0.247134
B    0.236708
C    0.187789
Name: assigned_room_type, dtype: float64
```

Among the top seven room types with the most bookings, the cancellation rates of room types A and G are higher than other room types, and the cancellation rate of room type A is as high as 44.5%.